

Integration of transcriptome and genome sequencing uncovers functional variation in human populations



Tuuli Lappalainen¹, M Sammeth^{2,3}, M Friedländer², PAC 't Hoen⁴, M Rivas⁵, J Monlong², M Gonzales-Porta⁶, N Kurbatova⁶, T Griebel^{2,3}, Ferreira P², M Barann⁷, T Wieland⁸, L Greger⁶, M van Iterson⁴, J Almløf⁹, P Ribeca³, I Pulyakhina⁴, D Esser⁷, E Lizano², M Sultan¹⁰, D MacArthur¹¹, I Padioleau¹, T Strom⁸, McCarthy M⁵, H Lehrach¹⁰, S Schreiber⁷, R Sudbrak¹⁰, A Carracedo Alvarez¹², S Antonarakis¹, Robert Häsler⁷, AC Syvanen⁹, GJ van Ommen⁴, A Brazma⁶, T Meitinger⁸, P Rosenstiel⁷, R Guigó², I Gut³, X Estivill², ET Dermitzakis¹, on behalf of the Geuvadis Consortium^{1,2,3,4,6,7,8,9,10,12}

1 Dept of Genetic Medicine and Development, University of Geneva; 2 Center for Genomic Regulation and UPF, Barcelona; 3 Centro Nacional de Análisis Genómico, Barcelona; 4 Center for Human and Clinical Genetics, Leiden University Medical Center; 5 Wellcome Trust Centre for Human Genetics, Oxford; 6 European Bioinformatics Institute, Hinxton; 7 Institute of Clinical Molecular Biology, University of Kiel; 8 Institute of Human Genetics, Helmholtz Zentrum München, Munich; 9 Department of Medical Sciences, Uppsala University; 10 Department of Vertebrate Genomics, Max Planck Institute for Molecular Genetics, Berlin; 11 Massachusetts General Hospital; 12 University of Santiago de Compostela

mRNA and miRNA sequencing of 465 samples from the 1000 Genomes project

Aims of the study: (1) How to do distributed RNA sequencing? (2) What can we learn of transcriptome variation and its genetic component by integrating genome and transcriptome data from hundreds of individuals? (3) Create one of the biggest reference datasets for transcriptomics.

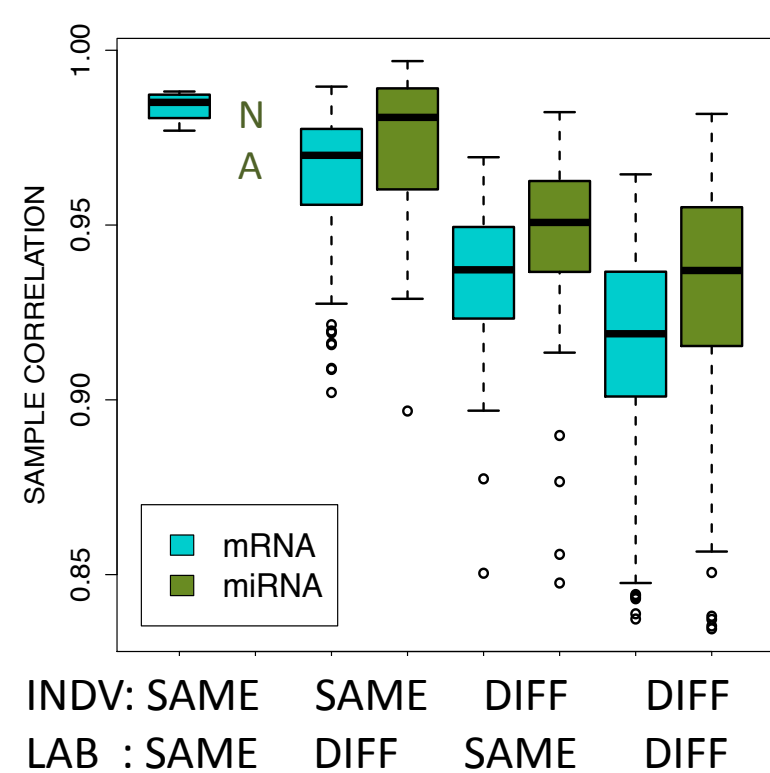
| | mRNA | miRNA |
|-----|------|-------|
| TSI | 93 | 89 |
| GBR | 94 | 94 |
| FIN | 95 | 93 |
| CEU | 91 | 87 |
| YRI | 89 | 89 |
| TOT | 462 | 452 |

RNA sequencing in 7 institutes with Illumina TruSeq protocol¹

- Random distribution of samples
- Replicates: 5 samples in each lab + 168 samples in two labs.
- Genotypes from 1000 Genomes²: 27 M total variants. 90% of samples in Phase1, the rest imputed from Omni2.5 M SNP data

| RNAseq reads | Individual median |
|---------------|-------------------|
| mRNA total | 58.4 M |
| mRNA QC pass | 48.8 M |
| miRNA total | 8.8 M |
| miRNA QC pass | 1.2 M |

| Quantifications | In >50% individuals |
|---------------------|---------------------|
| Genes | 14,779 |
| Exons | 141,951 |
| Transcripts | 74,533 |
| Splice junctions | 134,293 |
| Fusion genes | 5 |
| Transcribed repeats | 47,438 |
| RNA edited sites | 100 |
| miRNA genes | 706 |

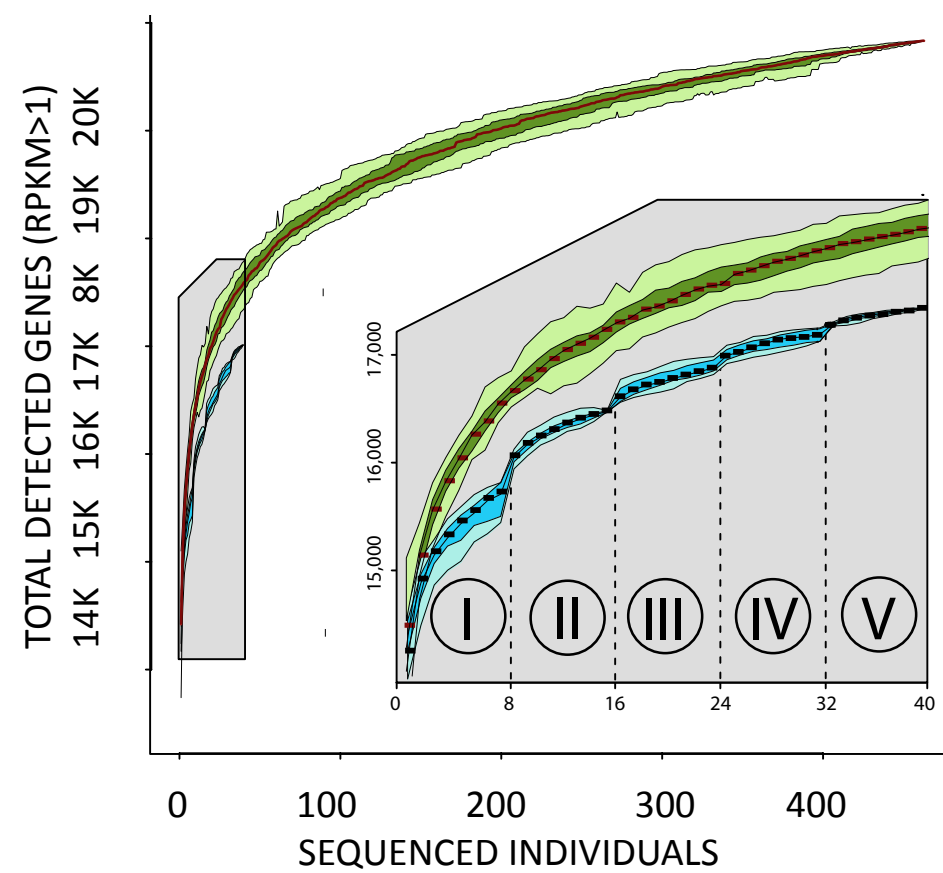


Distributed RNA-sequencing works well: technical variation due to laboratory effects is less than biological variation between samples

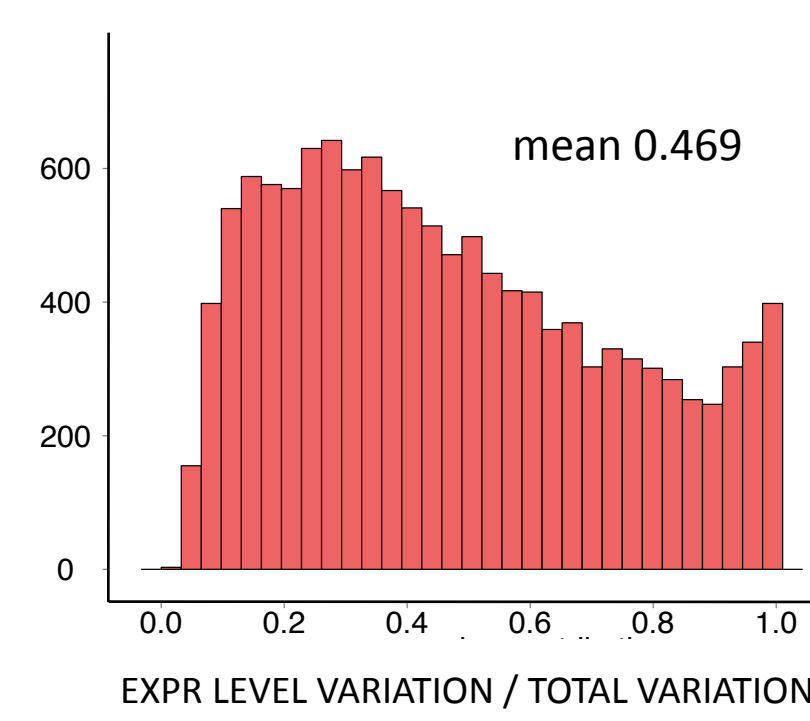
Data available : www.geuvadis.org ; ArrayExpress accessions E-GEUV-1, E-GEUV-2

Transcriptome variation within and between populations: mRNA, miRNA, and their interactions

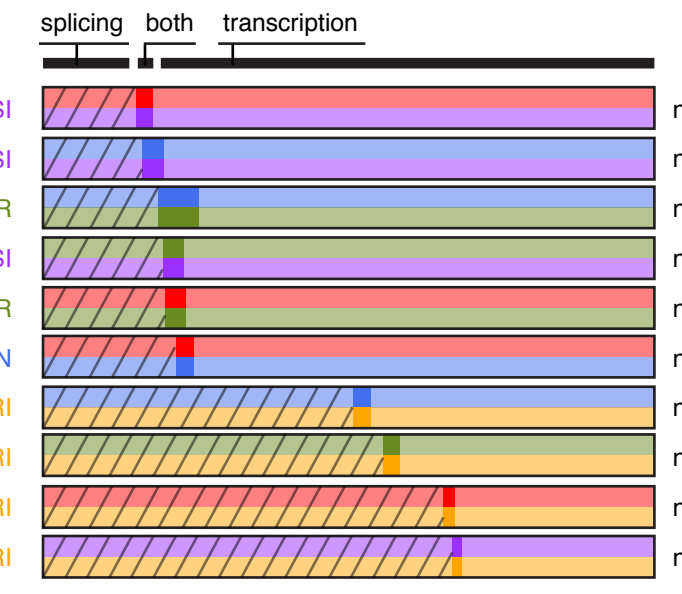
Population diversity adds 10% to gene detection



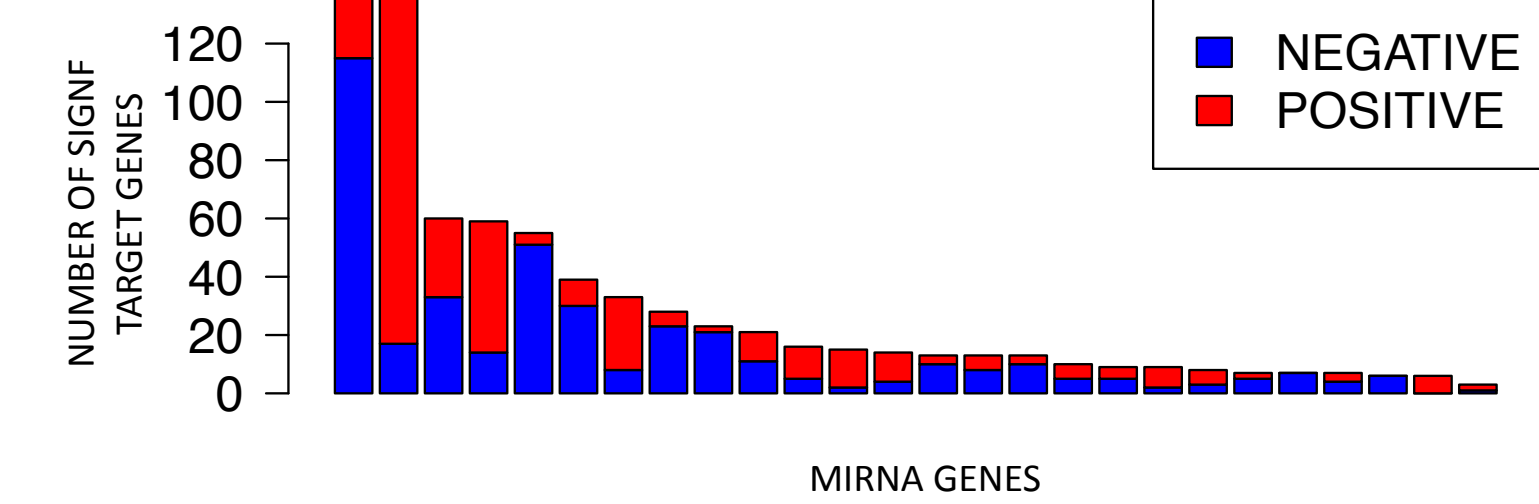
Gene expression levels and splicing contribute almost equally to transcription variation within populations.³



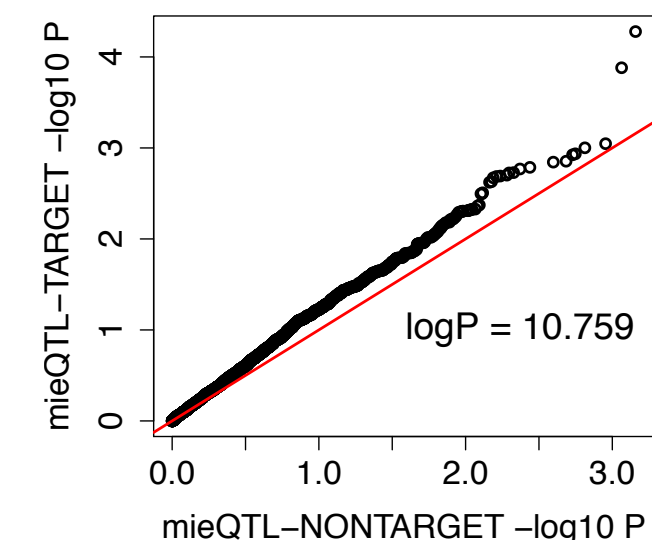
Differential splicing is more common between than within continents



The expression levels of 26 miRNAs correlate with predicted target⁴ quantifications in the population



miRNA cis-eQTLs have increased trans-eQTL signal with miRNA target genes

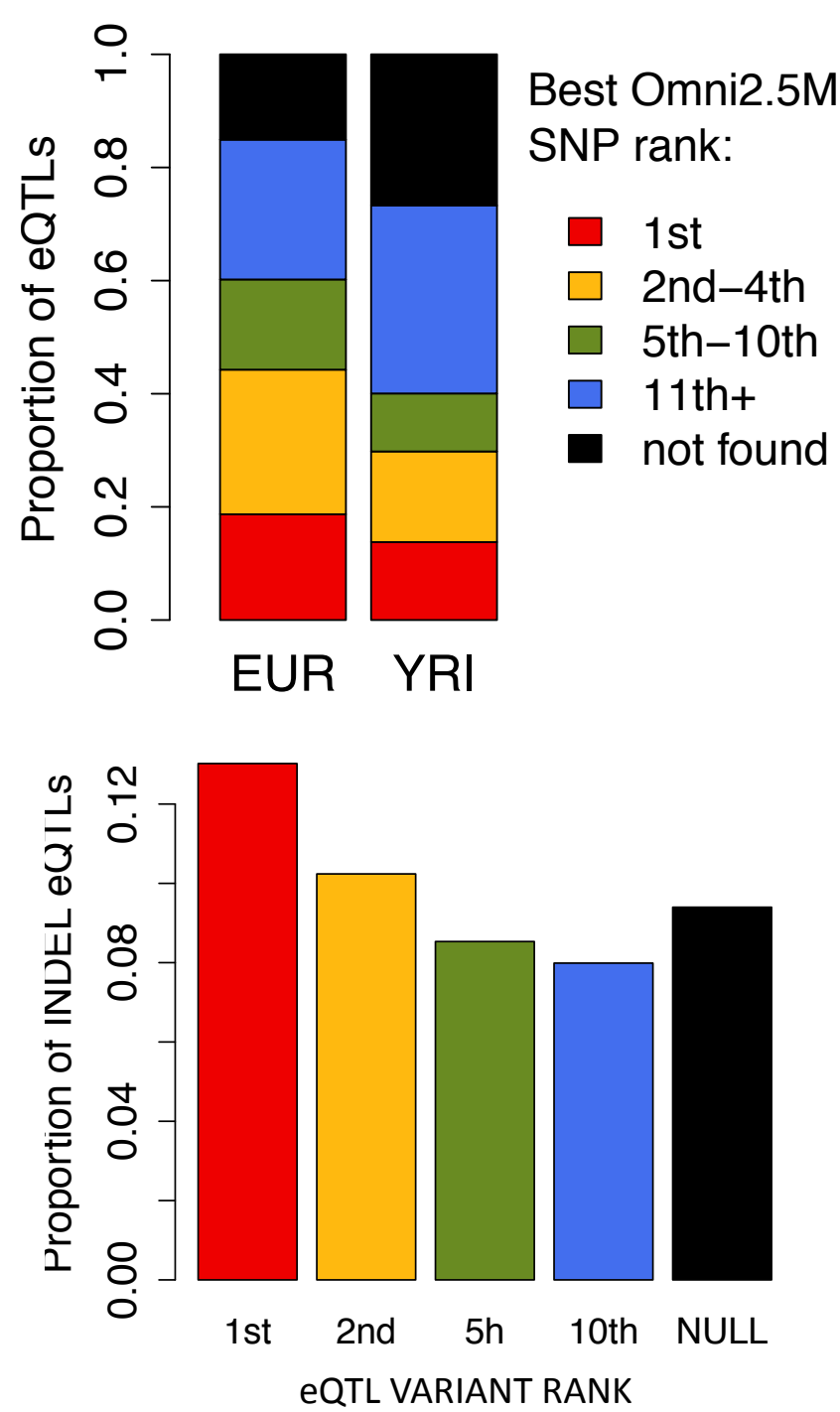


Thousands of expression and splicing cis-QTLs with increased discovery of causal variants

The majority of protein-coding genes and 10% of miRNA genes have a common regulatory variant.

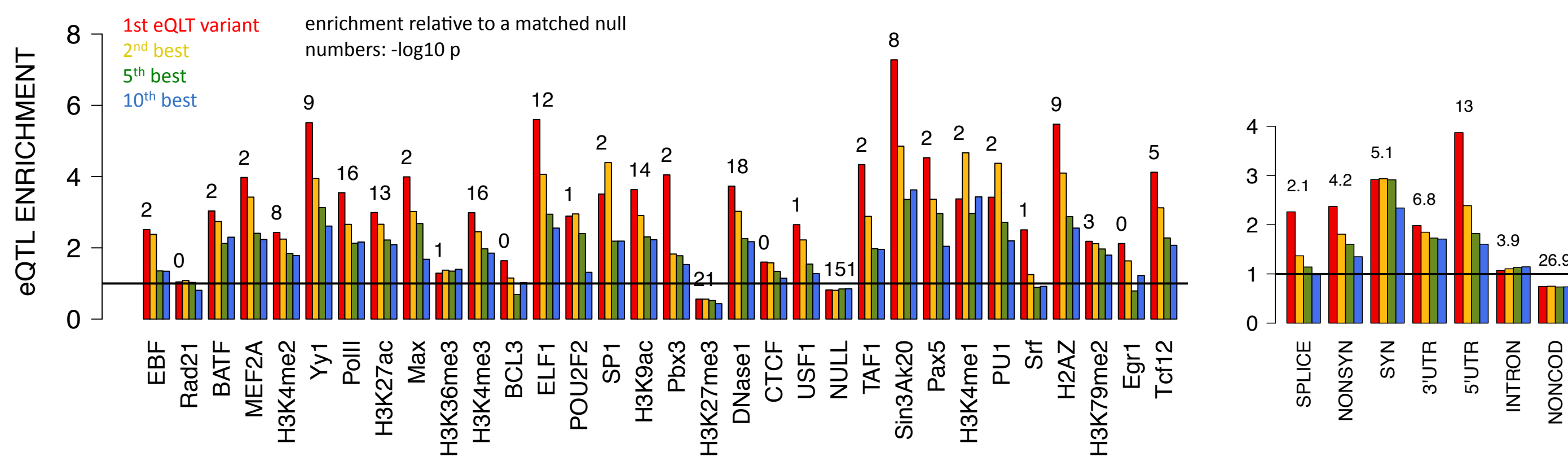
Genome sequencing data allows better discovery of the best-associating variants, and shows a significant enrichment of indels.

| | exon eQTLs (total 12982; FDR 5%) | asQTLs (total 16172; FDR 9.5%) | mi-eQTLs (total 644; FDR 5%) |
|-----------------|----------------------------------|--------------------------------|------------------------------|
| EUR (n=373) | 7486 | 2748 | 57 |
| YRI (n=89) | 2308 | 2511 | 15 |
| EUR & YRI union | 7877 | 4429 | 60 |

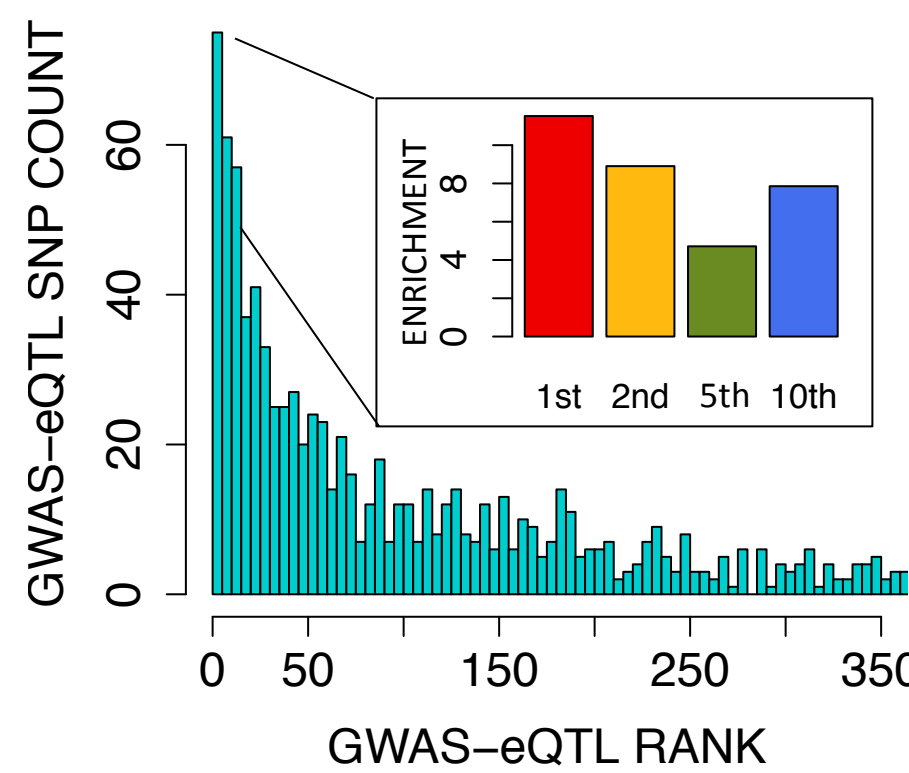


Enrichment of eQTLs in functional regions uncovers causes and effects of regulatory variation

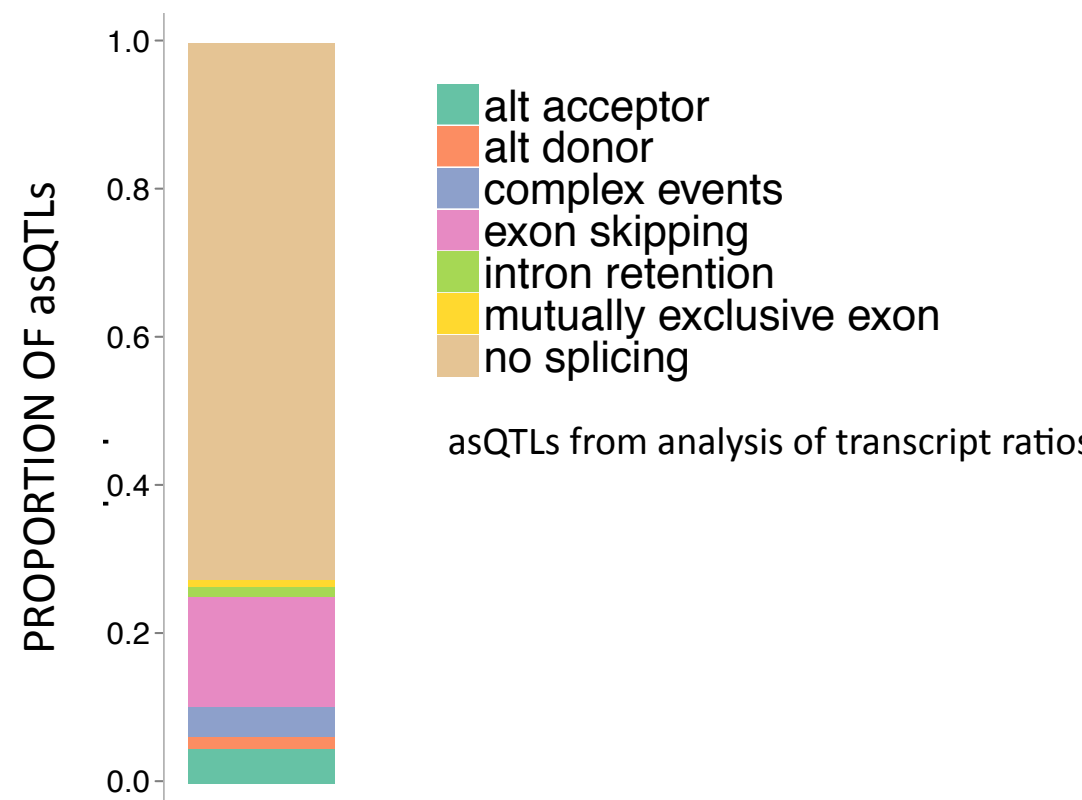
The best eQTL variants are significantly enriched in functionally annotated regulatory and coding regions (Ensembl Regulatory Build, Gencode v12), with an overrepresentation of especially promoter and enhancer annotations as well as splicing and nonsynonymous variants.



Variants from the NHGRI GWAS database are enriched among top eQTLs.

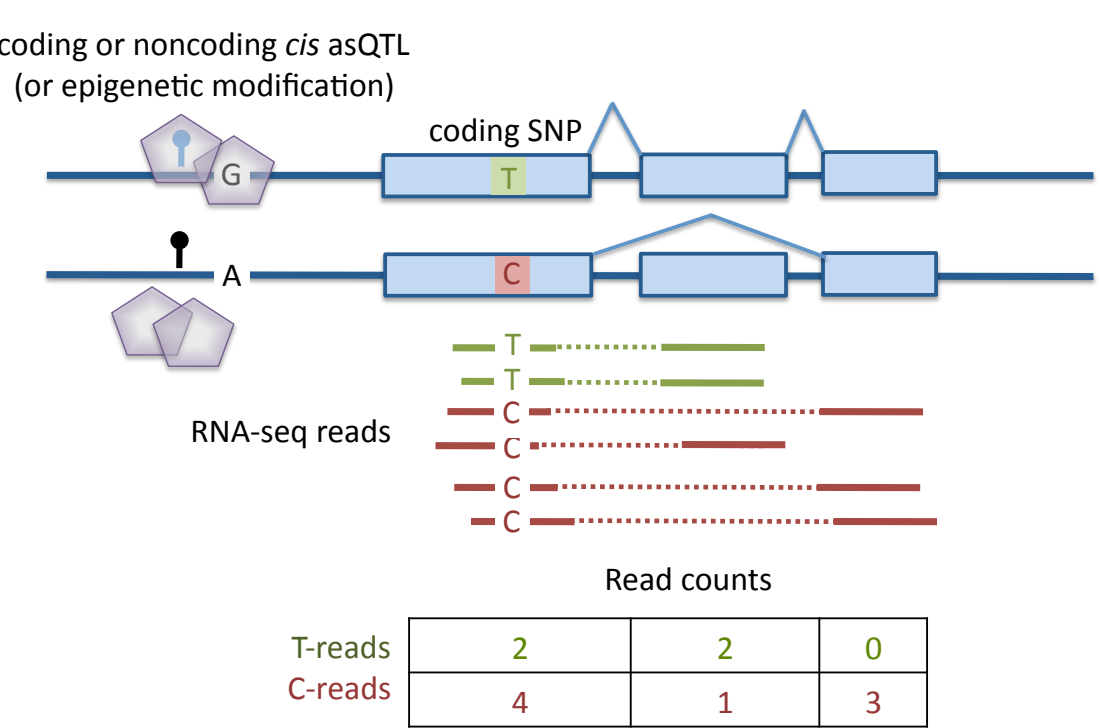


Most asQTLs affect UTR length rather than splicing.

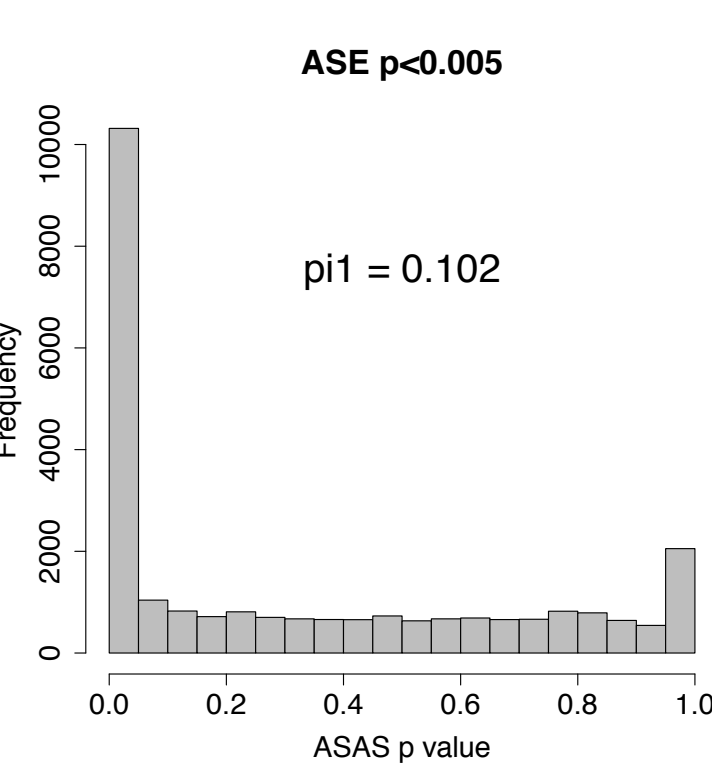


Variation in allelic expression is often driven by transcript structure variation and is dominated by rare effects

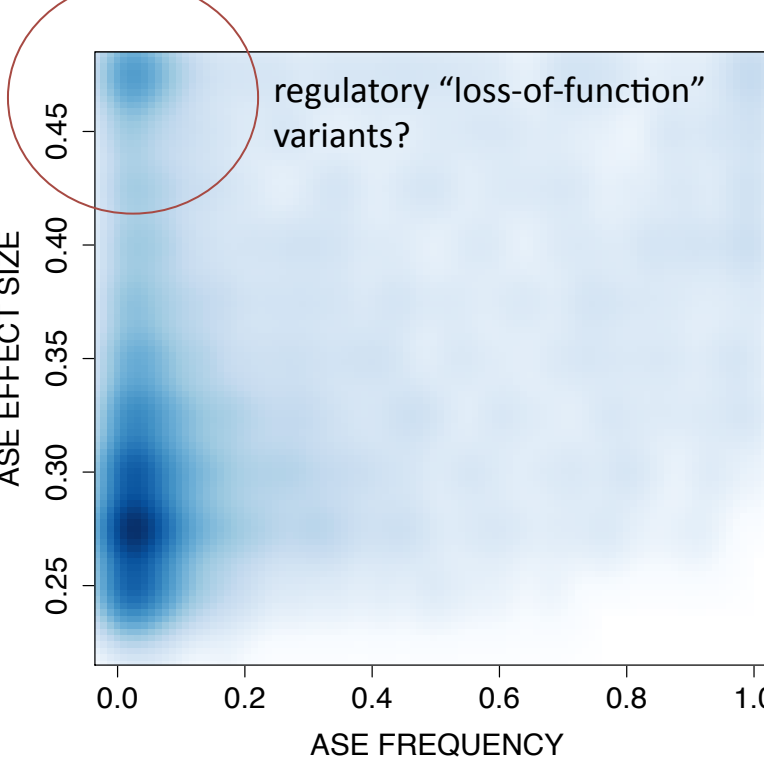
Detection of allele-specific expression (ASE) and transcript structure (ASTS)



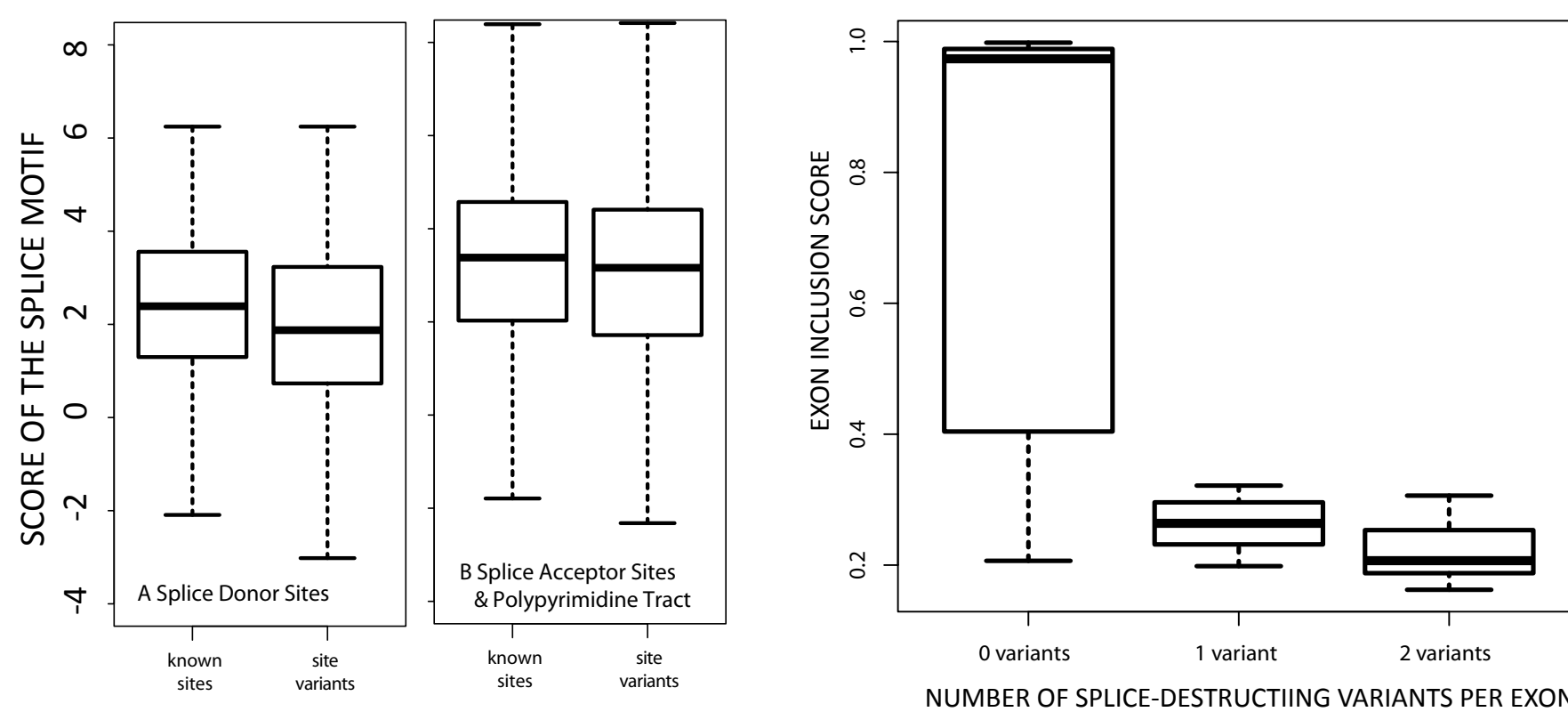
Much of ASE is driven by transcript structure changes



The majority of ASE is caused by rare events in the population



Functional characterization of loss-of-function variants



Genetic variants in splicing motifs significantly decrease both predicted and observed splicing efficiency.

Premature stop-codon variants lead to nonsense-mediated decay frequently but not always: better predictions needed

